

Coupling Tangent-Linear and Adjoint Models

Uwe Naumann¹ and Patrick Heimbach²

¹ Mathematics and Computer Science Division, Argonne National Laboratory
naumann@mcs.anl.gov, <http://www.mcs.anl.gov>

² Earth, Atmospheric and Planetary Sciences, Massachusetts Institute of Technology
heimbach@mit.edu, <http://mitgcm.org>

Abstract. We consider the solution of a (generalized) eigenvalue problem arising in physical oceanography that involves the evaluation of both the tangent-linear and adjoint versions of the underlying numerical model. Two different approaches are discussed. First, tangent-linear and adjoint models are generated by the software tool TAF and used separately. Second, the two models are combined into a single derivative model based on optimally preaccumulated local gradients of all scalar assignments. The coupled tangent-linear / adjoint model promises to be a good solution for small or medium sized problems. However, the simplicity of the example code at hand prevents us from observing considerable run time differences between the two approaches.

1 Introduction

Substantial effort in understanding the ocean circulation’s role in the variability of the climate system, on time scales of decades to millennia and beyond, is being directed at investigating the so-called ‘thermohaline’ circulation (THC). This refers to the contribution to the ocean circulation which is driven by density gradients and thus controlled by temperature and salinity properties and its associated fluxes (see e.g. [?], but see [?] for problems of its isolated discussion). It plays a crucial role in connecting the surface to the deep ocean through deep-water formation which occurs at some isolated convection sites at high latitudes mainly in the subpolar Atlantic ocean, such as the Labrador Sea and the Greenland-Irminger-Norwegian (GIN) Seas.

Because of its relatively small variability over the past few thousand years compared to potentially large changes on glacial-interglacial time scales, its dynamics may be well described by modal behavior of linear dynamics such as stochastically (or realistically) forced damped oscillations. Major changes in the THC are then ascribed to either stochastic forcing across the stability threshold, shifting the system between stable equilibria, or nonlinear effects such as self-sustained oscillations. In addressing what triggers instabilities in the climate system, comparatively little attention has been given to systems which are non-normal, i.e. which exhibit non-orthogonal eigenmodes (beginning with [?]; see [?] and references therein). Such systems, albeit linear, can undergo large transient amplifications. In [?] these ideas were applied to a simplified box model

of the THC. It serves as a paradigm model to capture some essential features of the circulation [?,?]. The authors demonstrated how, for a given optimality condition (norm), a set of initial conditions of temperature and salinity may be determined which maximize the transient amplification of the THC. From an automatic differentiation (AD) [?] point of view, their approach is interesting for two reasons: (i): their calculation involves both the tangent linear (TLM) and the adjoint (ADM) operator of the model, (ii): their approach can be generalized to a coupled atmosphere-ocean model of intermediate complexity, or even to a fully-fledged general circulation model (GCM), for which AD becomes a crucial ingredient in deriving the TLM and ADM.

For the present purpose, our focus remains on the simple box model, for which we wish to derive the TLM and ADM by means of AD. The box model is introduced in the following section. The optimality problem is stated in Section 3. Section 4 briefly describes how the TLM and ADM were generated by means of FastOpt’s AD tool TAF (Transformation of Algorithms in Fortran) [?,?]. The main reason for choosing this tool is that it has been applied successfully to generate efficient derivative code of MIT’s parallel ocean general circulation model (MITgcm), [?,?,?], and will thus play an important part in the generalization of the box model study to a GCM.

An algorithm for coupling the TLM and ADM is presented in Section 5. It is based on optimally preaccumulated gradients of scalar assignments [?,?]. Both approaches are compared and conclusions are drawn in Section 7 in the light of current work on a new infra-structure for the implementation of next-generation AD algorithms.

2 A simple model of the thermohaline circulation

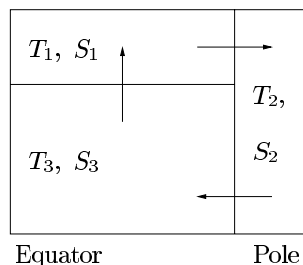


Fig. 1. A simple model of the thermohaline circulation

Some important aspects of linear dynamics of the THC in the North Atlantic are captured by a simple 3-box model [?,?], see Figure 1. Box 1 represents the shallow surface waters between the equator and the sub-polar ocean. Box 2 reflects the dynamics of the polar ocean and its role in deep water formation. Box 3 captures the deep ocean at low and mid latitudes. The model state is represented by properties of temperature T_i and salinity S_i in each of the three boxes, $i = 1, 2, 3$. The dynamics is governed by a set of equations; here, we only

state the sub-set for box 1 for the case $U \geq 0$, the full set of equations can be found in [?]:

$$\begin{aligned}
U &= u_0 \left\{ \rho_2 - \left[\frac{H_1}{H_3} \rho_1 + \left(1 - \frac{H_1}{H_3} \right) \rho_3 \right] \right\} \\
\rho_i &= -\alpha T_i + \beta S_i, \quad i = 1, 2, 3 \\
\frac{d}{dt} T_1 &= U(T_3 - T_1)/\mathcal{V}_1 + F_1^{relax} \\
\frac{d}{dt} S_1 &= U(S_3 - S_1)/\mathcal{V}_1 + F_1^{fw}
\end{aligned} \tag{1}$$

U denotes the net density-driven transport between the boxes with mean transport u_0 , density ρ_i is inferred via a linear equation of state with thermal and salinity expansion coefficients α and β , respectively, and tracer advection terms dT_1/dt , dS_1/dt for temperature and salinity. \mathcal{V}_i refers to the volume of the box i of height H_i ; F_1^{relax} , F_1^{fw} denote external forcing due to air-sea fluxes of heat (or temperature relaxation) and freshwater, respectively. The full set of parameter choices are given in [?]. The model is stepped forward in time by a standard upwind differencing scheme.

3 The optimality criterion

The time-evolving temperature and salinity in each box together define a six-dimensional state vector \mathbf{P} , whose perturbation we denote as

$$\delta \mathbf{P} = [\delta T_1, \delta T_2, \delta T_3, \delta S_1, \delta S_2, \delta S_3]^T, \tag{2}$$

For a given model (or function) F , its tangent linear operator (or Jacobian) $F' = \mathbf{dF}(\tau, 0)$ defines a mapping of the model's initial ($t = 0$) to its final ($t = \tau$) tangent space,

$$\delta \mathbf{P}(\tau) = \mathbf{dF}(\tau, 0) \delta \mathbf{P}(0) \tag{3}$$

In the tangent spaces at times t we introduce scalar products $\langle \mathbf{x}, \mathbf{y} \rangle_t = \mathbf{x} \mathbf{M}_t \mathbf{y}$, with matrix representations \mathbf{M}_t , and their induced norms $\|\cdot\|_t$, which may be identified with the strength of the THC transport U [?],

$$U^2(t) = \|\delta \mathbf{P}\|_t^2 = \langle \delta \mathbf{P}(t), \delta \mathbf{P}(t) \rangle = \delta \mathbf{P}^T(t) \mathbf{M}_t \delta \mathbf{P}(t) \tag{4}$$

With respect to these norms we can define amplification rates of initial perturbations over an integration period τ

$$a = \frac{\|\delta \mathbf{P}(\tau)\|_\tau}{\|\delta \mathbf{P}(0)\|_0} \tag{5}$$

We then wish to find those initial perturbations $\delta \mathbf{P}_0$ which maximize the amplification rate over the integration period τ , i.e. whose initial norm is unity and final norm is maximal. This constrained maximization problem can be expressed

as an unconstrained problem by means of the Lagrange multiplier method, i.e. we seek

$$\max_{\delta \mathbf{P}_0} \left\{ \langle \delta \mathbf{P}(\tau), \delta \mathbf{P}(\tau) \rangle + \lambda [\langle \delta \mathbf{P}(0), \delta \mathbf{P}(0) \rangle - 1] \right\} \quad (6)$$

Differentiating w.r.t. $\delta \mathbf{P}_0$ to infer the extrema and using the definition of the adjoint of the tangent linear operator \mathbf{dF} yields a generalized eigenvalue problem of the form

$$\mathbf{dF}^T \mathbf{M}_\tau \mathbf{dF} \mathbf{e} = \lambda \mathbf{M}_0 \mathbf{e} \quad (7)$$

with (generalized) eigenvalues and vectors λ and \mathbf{e} , respectively. Here, \mathbf{dF}^T denotes the adjoint (transpose) operator.

This problem can be readily solved in MATLAB [?]. However, in the present context of an implementation in FORTRAN, and in view of future extensions to a fully-fledged ocean model, an alternative iterative strategy is adopted: First, at a given iteration n , a given estimate \mathbf{e}_n is used to compute

$$\mathbf{y}_n = \mathbf{dF}^T(\tau, 0) \mathbf{M}_\tau \mathbf{dF}(\tau, 0) \mathbf{e}_n \quad (8)$$

by means of the tangent linear and adjoint operator. Second, \mathbf{y}_n and \mathbf{e}_n are provided to an Implicitly Restarted Arnoldi Iteration Routine, ARPACK [?], which yields new estimates \mathbf{e}_{n+1} and λ_{n+1} . This method enables an efficient computation of the leading eigenvectors and eigenvalues for large-scale applications, for which the dimension of the state space is large ($\sim 10^7$ for ocean GCM's).

4 Automatic TLM and ADM Generation via TAF

The generation of the tangent-linear model via TAF is straightforward: providing the dependent variables (final state) and independent variables (initial state), TAF produces readable tangent linear code which can be readily used.

For the adjoint code, a few interventions were made in order to deal with state variables that are required in evaluating the derivative expressions in reverse order. The time stepping loop was split into a two-level checkpointing (see e.g. [?,?,?]) according to `nTimeSteps = nOuter * nInner`. The model state is then stored to disk once every `nOuter` time steps. At these instances the model picks up to recompute over an interval of `nInner` time steps. Over this interval only it stores the required variables at each time step to common blocks to be available in reverse mode. This enables an efficient balance of storing vs. recomputation at the cost of one additional model integration and some extra memory.

5 Coupling TLM and ADM

The model implements a vector function $F : \mathbb{R}^6 \rightarrow \mathbb{R}^6$ such that $\mathbf{y} = F(\mathbf{x})$, where \mathbf{x} represents the initial and \mathbf{y} the final state, as in Section 4. The Jacobian at some point \mathbf{x}_0 is denoted by $F' \in \mathbb{R}^{6 \times 6}$. At each iteration of the algorithm

used in Section 4 we compute $\dot{\mathbf{y}} = F' \cdot \dot{\mathbf{x}}$ followed by $\bar{\mathbf{x}} = (F')^T \cdot \dot{\mathbf{y}}$. Instead of generating two separate tangent-linear and adjoint versions of F , we now combine both into one derivative model.

An augmented version of TLM performs the following steps for all assignments $v = f(\mathbf{u})$ in the original model: (A) generation of code list; (B) local activity analysis and linearization; (C) optimal preaccumulation of gradient using the algorithm described in [?]; (D) computation of the inner product $\dot{v} = f' \cdot \dot{\mathbf{u}}$, where f' denotes the gradient of f ; and (E) storage of entries of f' on a stack (also known as the *tape*). An example is discussed in Section 6.

The local gradients are restored in reverse order to propagate adjoints backward through the code for each statement $v = f(\mathbf{u})$ as $\bar{\mathbf{u}} = \bar{v} \cdot f'$. In contrast with the method described in Section 4, these gradients are computed only once and used in both TLM and ADM. A downside of this approach is the requirement to store the entire derivative information. This may not be feasible for larger problems. On the other hand, local gradients of scalar assignments can be preaccumulated optimally [?,?], a procedure that may lead to a further reduction of the computational effort.

Steps (A)–(D) are fully automated at the graph level and implemented as part of the XAIFBooster library that is being developed at Argonne National Laboratory. It uses a language-independent intermediate format called XAIF [?] for the semantic transformation of numerical programs. A protocol of how the transformation of the assignments in subroutine `box_timestep` was done can be found under <http://www-unix.mcs.anl.gov/~naumann/iccsa03>. The actual code generation was done manually. Work is in progress to interface the library with front-ends for both C/C++ and Fortran 95 in order to provide a new software tool for AD.

From the viewpoint of differentiation the information provided by the output of the XAIFBooster module is nearly sufficient to generate both tangent-linear and adjoint models in general. We have correct (and efficient) gradient code for all statements and the information on the corresponding locally dependent and independent variables. All we need in addition are correct unparsers and, for adjoint models, a mechanism to reverse the flow of control. The latter turns out to be trivial for the simple box model at hand. Although this approach is certainly feasible, it may not be optimal in many cases.

Full automation of the entire process is the subject of ongoing work at Argonne National Laboratory and Rice University. In collaboration with scientists at MIT a new infrastructure for the implementation of AD algorithms is being developed. This work is supported by the National Science Foundation under its Information Technology Research program. See <http://www.autodiff.org/ACTS> for further information on the project.

The source code of the coupled tangent-linear/adjoint model that is based on optimally preaccumulated local gradients of scalar assignments can be accessed online at <http://www-unix.mcs.anl.gov/~naumann/iccsa03>.

6 Optimal Statement-Level Gradient Accumulation

Using the first assignment from the subroutine `box_timestep`, we sketch the algorithm used to preaccumulate the local gradients. A more complete description of the procedure as well as proofs for its optimality regarding the computational effort can be found in [?].

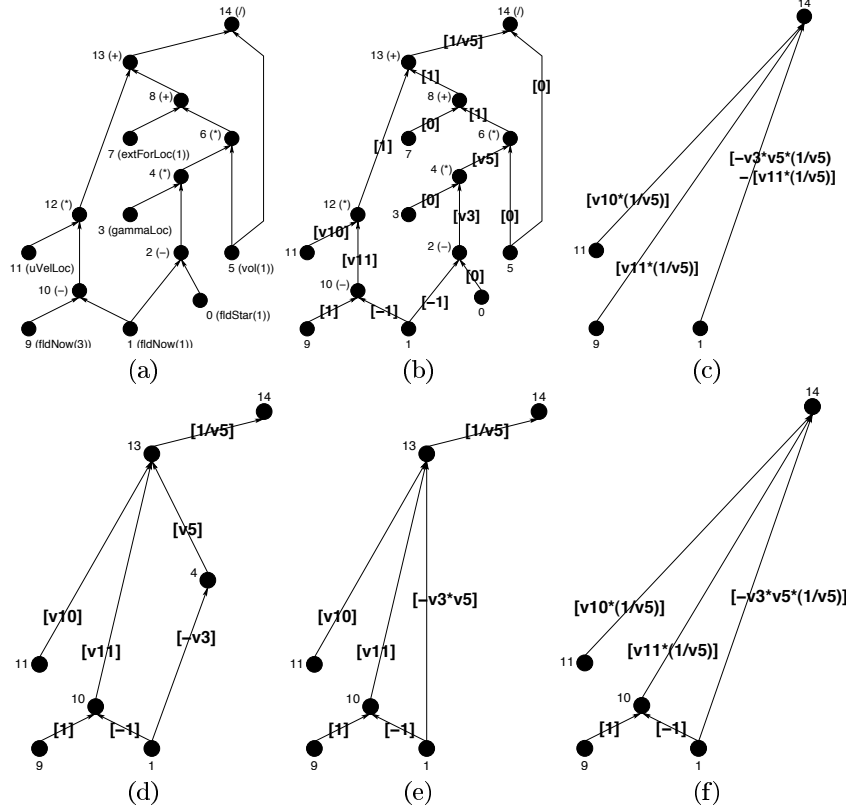


Fig. 2. Manipulation of the computational graph

Original Statement. The original statement has the following form.

```
dFldDt(1) = ( extForLoc(1)
&   + gammaLoc*( fldStar(1) - fldNow(1) )*vol(1)
&   + uVelLoc*( fldNow(3)-fldNow(1))) / vol(1)
```

Its computational graph G is shown in Figure 2(a). We are interested in the gradient of the dependent variable $dFldDt(1)$ with respect to the independent variables $fldNow(1)$, $fldNow(3)$, and $uVelLoc$. Our objective is to convert G into the bipartite graph in Figure 2(c) whose edge labels represent exactly the entries of the gradient [?].

Code List. The aim of an optimized code list is to assign to auxiliary variables only those intermediate values that are used by the computation of some local partial derivative. Here this applies only to the value of `fldNow(3)-fldNow(1)`. Thus, the number of new variables to be generated in the augmented TLM can be kept minimal. The following two assignments are generated.

```
xb_aux_4=fldNow(3)-fldNow(1)
dFldDt(1)=( extForLoc(1) + gammaLoc*(fldStar(1) - fldNow(1))
&    *vol(1)+uVelLoc* xb_aux_4) / vol(1)
```

Local Activity Analysis. The statement has three *active* [?] inputs (`fldNow(1)`, `fldNow(3)`, and `uVelLoc`) represented by the independent vertices 1, 9, and 11 in **G**. Vertices that can be reached from some independent vertex via a path in **G** are active (2,4,6,8,10,12,13,14). All other vertices are passive (0,3,5,7). Edges emanating from passive vertices are labeled with 0 in Figure 2(b), as the local partial derivative with respect to a passive variable vanishes identically.

Local Partial Derivatives and Elimination Procedure. The computational graph is linearized by attaching the local partial derivatives of the elemental functions with respect to their arguments to the corresponding edges. The result is shown in Figure 2(b), where the partial derivatives are enclosed in square brackets.

Constant labels of incident edges can be *folded* [?] at compile-time. If (i, j) is labeled with $c_{j,i}$ and (j, k) with $c_{k,j}$ and both $c_{j,i}$ and $c_{k,j}$ are constants, then the value of $c_{j,i} \cdot c_{k,j}$ can be evaluated and (i, j) can be *front eliminated* [?] at compile time. Notice that (j, k) must be the only edge emanating from j because **G** is a single-expression-use graph as in [?].

Trivial edges, labeled with 1 or -1, are *back eliminated* at no extra cost. This procedure leads to the graph in Figure 2(d). The front elimination of either edges (1, 10) or (9, 10) would not preserve optimality because vertex 10 has more than one predecessor and its successor is not equal to the dependent vertex. Refer to [?] for proofs of these results and a more complete description of the constant folding algorithm.

The graph in Figure 2(d) is reduced to the bipartite graph in Figure 2(c), which represents the local gradient by the optimal preaccumulation algorithm for single-expression-use graphs described in [?]. Vertex 4 is eliminated (leading to the graph in Figure 2(e)) followed by 13 (Figure 2(f)) and 10. The latter does not result in any scalar floating-point multiplications because all edges leading into it are trivial. As a common subexpression in all gradient entries $1/v5$ is assigned to an auxiliary variable. The code resulting from this elimination is shown below.

```
c_14_13=1./vol(1);          c_14_11=xb_aux_4*c_14_13
c_14_10=uVelLoc*c_14_13;    c_14_1=-gammaLoc*vol(1)*c_14_13
c_14_9=c_14_10;            c_14_1=c_14_1-c_14_10
```

Computation of Directional Derivative. Finally, the entries of the gradient are used to compute the directional derivative $dF_{ldDt_d}(1)$ as

$$dF_{ldDt_d}(1) = c_{14_11} * uVelLoc_d + c_{14_1} * fldNow_d(1) + c_{14_9} * fldNow_d(3).$$

Following this, the local gradient is stored on the tape.

Computation of Adjoints. After the values for c_{14_1} , c_{14_9} , and c_{14_11} are restored, they are used to compute the corresponding adjoints as

$$\begin{aligned} fldNow_a(3) &= fldNow_a(3) + c_{14_9} * dF_{ldDt_a}(1) \\ fldNow_a(1) &= fldNow_a(1) + c_{14_1} * dF_{ldDt_a}(1) \\ uVelLoc_a &= uVelLoc_a + c_{14_11} * dF_{ldDt_a}(1); \quad dF_{ldDt_a}(1) = 0 \end{aligned}$$

A formal proof for the statement-level optimality of the above approach can be found in [?].

7 Comparison and Conclusion

Table 1. Runtime (in sec) of 10,000 TLM/ADM evaluations

| | Separate TLM and ADM | Coupled TLM and ADM |
|-------------------------------|----------------------|---------------------|
| Without compiler optimization | 286 | 265 |
| With compiler optimization | 167 | 183 |

We compared the elapsed times for 10,000 evaluations of the coupled TLM/ADM. The results are displayed in Table 1. They depend strongly on the code optimizations performed by the compiler (g77, in this case, with or without the `-O3` option activated). The coupled TLM/ADM approach generates more efficient source code and is superior if compiler optimization is switched off. On the other hand, its memory requirements are higher than those of the adjoint code generated by TAF and featuring a two-level checkpointing strategy. The repeated recomputations performed by the TAF-generated adjoint code do not carry much weight in our very simple example. Moreover, the real power of the local gradient accumulation algorithms proposed in [?,?] cannot be exploited either, because of the relative simplicity of the single statements. Enabling full compiler optimization turns the emphasis with regard to the overall performance toward the memory traffic. There the coupled TLM/ADM loses because of the requirement to store all the local gradients of the assignments. The number of floating-point operations to be performed during the accumulation of these gradients is small.

The theoretical optimality of the local gradient accumulation routine is not sufficient to compensate for the runtime increase from higher memory requirements. This result is not surprising because the storage of local gradient entries increases the amount of memory required by a factor that is equal to the number of active arguments on the right-hand side of the assignment, compared with the strategy of storing values of intermediate variables. If, however, the preaccumulation can be applied to local Jacobians at the basic-block level, this is no longer

the case. There the minimization of the number of edges in the local linearized computational graph can lead to a decrease in the amount of memory required.

We conclude that the development of efficient techniques for coupling tangent-linear and adjoint models represents a challenging research area in the field of automatic differentiation. Most likely, useful approaches will represent heuristic compromises between local preaccumulation techniques and hierarchical checkpointing algorithms. The ideas presented in this paper should be applied to more complex examples allowing for a better exploitation of the optimizations performed by the local preaccumulation algorithm. An appropriate source transformation platform for automatic differentiation is currently being developed as part of the ACTS project.

Acknowledgments

This work was supported by the National Science Foundation's Information Technology Research Program under Contract OCE-0205590 ("Adjoint Compiler Technology and Standards").

Naumann was supported by the Mathematical, Information, and Computational Sciences Division subprogram of the Office of Advanced Scientific Computing Research, U.S. Department of Energy, under Contract W-31-109-ENG-38.

The adjoint code of Section 4 was generated via TAF (Transformation of Algorithms in Fortran) developed by FastOpt Inc., Hamburg, Germany. A matlab version of the original box model was kindly provided by E. Tziperman. His invitation and funding by the Weizmann Institute of P.H. to initiate the work on optimal perturbations are gratefully acknowledged. P.H. is supported in part by the ECCO (Estimating the Circulation and Climate of the Ocean) Consortium with funding from NSF, NASA and ONR.